



# New improvements in the use of dependence measures for sensitivity analysis and screening

Matthias de Lozzo, Amandine Marrel

## ► To cite this version:

Matthias de Lozzo, Amandine Marrel. New improvements in the use of dependence measures for sensitivity analysis and screening. *Journal of Statistical Computation and Simulation*, 2016, 86 (15), pp.3038-3058. hal-01090475v2

**HAL Id: hal-01090475**

**<https://hal.science/hal-01090475v2>**

Submitted on 4 Dec 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# New improvements in the use of dependence measures for sensitivity analysis and screening

Matthias De Lozzo\*      Amandine Marrel†  
CEA, DEN, DER, F-13108 Saint Paul Lez Durance, France

December 4, 2015

## Abstract

Physical phenomena are commonly modeled by time consuming numerical simulators, function of many uncertain parameters whose influences can be measured via a global sensitivity analysis. The usual variance-based indices require too many simulations, especially as the inputs are numerous. To address this limitation, we consider recent advances in dependence measures, focusing on the distance correlation and the Hilbert-Schmidt independence criterion. We study and use these indices for a screening purpose.

Numerical tests reveal differences between variance-based indices and dependence measures. Then, two approaches are proposed to use the latter for a screening purpose. The first approach uses independence tests, with existing asymptotic versions and spectral extensions; bootstrap versions are also proposed. The second considers a linear model with dependence measures, coupled to a bootstrap selection method or a Lasso penalization. Numerical experiments show their potential in presence of many non-influential inputs and give successful results for a nuclear reliability application.

**Keywords:** sensitivity analysis, screening, dependence measures, independence tests, bootstrap, HSIC.

## 1 Introduction

Numerical simulators are widely used in the industry for the representation of physical phenomena (Santner et al, 2003). Such models take as input a high number of numerical and physical explanatory variables. The information on these underlying input parameters is often limited or uncertain. Commonly, the uncertainties on the input parameters are modeled by probabilistic distributions. Then, the objective is to assess how these uncertainties can affect the model output. For this, computer experiments methodologies based on statistical advanced techniques are useful (de Rocquigny et al, 2008).

Sensitivity Analysis (SA) methods allow to answer the question “How do the input parameters variations contribute, qualitatively or quantitatively, to the variation of the output?” (Saltelli et al, 2008). More precisely, these tools can detect non-significant input parameters in a screening context, determinate the most significant ones, measure their respective contributions to the output or identify an interaction between several inputs which impacts strongly the model output. In such a way, engineers can guide the characterization of the model by reducing the output uncertainty: they can calibrate the most influential inputs and fix the non-influential ones to nominal values. Many surveys on SA exist in the literature, such as (Kleijnen, 1997), (Frey and Patil, 2002) or (Helton et al, 2006); they divide the SA into two sub-domains: the Local Sensitivity Analysis (LSA) and the Global Sensitivity Analysis (GSA). The first one studies the effects of small input perturbations around nominal values on the model output. Usually this deterministic approach considers the partial derivatives of the model at a specific value of the input vector (Cacuci, 1981). The second sub-domain of SA considers the impact of the input uncertainty on the output over the whole variation domain of uncertain inputs, that is why it is called Global SA (Saltelli et al, 2008).

The GSA can be used for quantitative or qualitative purposes, specific tools being dedicated to each aim. From one hand, quantitative GSA methods supply an order of the input parameters which is function of their dependence to the output. Among them, the Derivative-based Global Sensitivity Measures (DGSM) consider the mean of the model gradient over the whole input domain (Lamboni et al, 2013), not at a specific point like in LSA (Cacuci, 1981). Another approach is based on the decomposition of the output variance; in particular, the Sobol’ indices are widely used and measure the proportion of the output variance explained by each input parameter (Sobol, 1993). Other authors propose to consider all the probabilistic distribution and not only the variance, comparing the distribution of the output conditioned by an input parameter with the unconditioned one (Borgonovo, 2007).

---

\*matthias.delozzo@cea.fr

†amandine.marrel@cea.fr

From the other hand, qualitative GSA uses less costly tools coming from the screening field. These methods can detect the input-output dependences and separate the input parameters into two groups: the non-significant ones and the significant ones. Despite of the criticisms with respect to the underlying hypotheses (Saltelli and Annoni, 2010), the basic screening tool is the one-at-a-time (OAT) design which consists in changing the values of each input parameter in turn from a control level scenario to a lower or upper level and measuring the evolution magnitude of the output (Daniel, 1958). Another method is the Morris design which consists in the repetition of many OAT designs, in order to get a mean value and a standard deviation for each input elementary effect (Morris, 1991). Other screening methods are currently used, such as the sequential bifurcation in a sparse context, when the number of significant input parameters is considerably lower than the total one which is greater than the number of observations (Bettonvil and Kleijnen, 1997). When the number of observations and the number of input parameters are of the same order, factorial fractional designs and other popular designs of experiments can be applied (Montgomery, 2006). Very recently, the use of Sobol' indices for sparse problems has been investigated (De Castro and Janon, 2014), in a screening framework where the effective dimension is much lower than the number of input parameters (Caflich et al, 1997).

Among all these GSA methods, the quantitative ones like Sobol' indices give a more accurate information about the dependence between the input parameters and the model output, while the qualitative methods are more imprecise. Moreover, Sobol' indices have been applied to many industrial problems in order to reduce the output variance. Nevertheless, these methods require many thousands of computer experiments in order to build reliable estimators of the sensitivity indices. Moreover, the number of required simulations is proportional to the number of inputs so as to preserve the precision of the sensitivity index estimator. Consequently, in the presence of a costly numerical simulator, quantitative GSA can not be performed directly, for high-dimensional problems.

A first alternative consists in replacing the computer code by a surrogate model and computing a quantitative GSA on this model. For example, Marrel et al (2009) and Sudret (2008) estimate the Sobol' indices thanks to Gaussian process models and polynomial chaos expansions respectively. However, the estimator accuracy depends on the precision of the surrogate model which can be weak if the learning sample is not enough representative. Moreover, the construction of the surrogate model in a high dimensional context (several decades of input parameters) is still an open problem.

Another alternative consists in using cheaper sensitivity indices which are potentially less accurate than the Sobol' ones but easier to compute (smaller CPU time). Qualitative GSA methods previously cited are commonly used to this aim. Nevertheless they often require either strong hypotheses on the model such as linearity, monotony or absence of interactions, or a number of observations much greater than the number of input parameters. The non-respect of these assumptions can lead to incorrect quantitative conclusions. Moreover, many of these screening methods consider specific design of experiments which can not be reused for other studies. Recently, new dependence measures removing these limitations have been developed by statisticians (Gretton et al, 2005; Székely et al, 2007) and applied in genomics, imagery or cross-language information retrieval (Blaschko et al, 2013). They have been studied in the field of global sensitivity analysis: they seem more robust than Sobol' indices, promising in a screening aim and can provide an information complementary to the Sobol' indices (Da Veiga, 2014). They can also make easier the metamodel construction by reducing the input number or guiding it in a sequential way; then, a quantitative GSA is performed to obtain an information more accurate on the input parameters identified as significant by the qualitative GSA.

**In this paper, we focus our attention on the use of these new dependence measures for qualitative GSA. We propose a guidance to use several independence tests based on these measures for a screening purpose: ones based on the estimator of the sensitivity index directly, others based on a linear decomposition and model selection methods. Some of these tests appear in literature and other ones are developed here. We also performed different numerical experiments to study the behavior of the dependence measures and compare the different proposed tests.**

Firstly, we present in Section 2 some dependence measures for the sensitivity of an output with respect to an input parameter. Secondly in Section 3, we deal with asymptotic and non-asymptotic statistical tests based on these dependence measures for feature selection. In Section 4, we propose a linear model associated to these dependence measures and build bootstrap tests and penalized regression techniques. Numerical experiments are tested on analytical models in Section 5, starting with a questioning around the meaning and the complementarity of the different sensitivity indices: "What sensitivity indices to what situation?". Finally, the significance tests based on HSIC are applied in Section 6 to a nuclear reliability application.

## 2 Sensitivity indices based on dependence measures

We consider a computer code  $Y = f(X_1, \dots, X_d)$  whose output  $Y$  and input parameters  $X_1, \dots, X_d$  belong to some measurable spaces  $\mathcal{Y}$ ,  $\mathcal{X}_1, \dots, \mathcal{X}_d$ . We note  $X = (X_1, \dots, X_d)$  the input vector of  $f$ .  $\mathcal{Y}$  and  $\mathcal{X}_k$  are commonly

equal to  $\mathbb{R}$ , for all  $k \in \{1, \dots, d\}$ , but sometimes engineers are in front of more complex situations where  $X_k$  or  $Y$  can be a vector, a time- or a space-discretized function, and so on. Because of their uncertainty (lack of knowledge, measuring accuracy, ...), the  $d$  input parameters are considered as realizations of random variables whose laws are perfectly known. Consequently, the output  $Y$  is also a random variable whose probability distribution is usually unknown and unapproachable because of the curse of dimensionality. **Sensitivity analysis aims at detecting and measuring the impact of each input uncertainty on the computer code output.**

We present in the following some measures of the dependence between an input parameter  $X_k$  and the output  $Y$  of the model  $f$ . The associated estimators are built using  $\left(X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)}\right)_{1 \leq i \leq n}$ , a  $n$ -sample of  $(X_1, \dots, X_d, Y)$ .

## 2.1 The Pearson's and Spearman's correlation coefficients

First of all, we can cite naive importance measures such as Pearson's and Spearman's correlation coefficients (see, e.g., Kendall and Stuart, 1977). These quantities evolve in the interval  $[-1, 1]$ , reaching the bounds for a total correlation between the variables  $X_k \in \mathcal{X}_k \subset \mathbb{R}$  and  $Y \in \mathcal{Y} \subset \mathbb{R}$  and equaling zero for an absolute uncorrelation. The Pearson's one is defined by  $\rho(X_k, Y) = \text{Cov}(X_k, Y) / \sqrt{\mathbb{V}[X_k] \mathbb{V}[Y]}$  and is estimated by  $\rho_n(X_k, Y) = \frac{\sum_{i=1}^n (X_k^{(i)} - \bar{X}_k)(Y^{(i)} - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_k^{(i)} - \bar{X}_k)^2 \sum_{i=1}^n (Y^{(i)} - \bar{Y})^2}}$  where  $\bar{X}_k = \frac{\sum_{i=1}^n X_k^{(i)}}{n}$  and  $\bar{Y} = \frac{\sum_{i=1}^n Y^{(i)}}{n}$ . The Spearman's correlation coefficient  $\tilde{\rho}_n(X_k, Y) = 1 - \frac{6 \sum_{i=1}^n d_{ik}^2}{n(n^2 - 1)}$  is a version of the Pearson's one applied on the ranks of  $\left(X_k^{(i)}, Y^{(i)}\right)_{1 \leq i \leq n}$ , where  $d_{ik} = \text{rank}\left(X_k^{(i)}\right) - \text{rank}\left(Y^{(i)}\right)$ . Asymptotically, the associated statistics  $t_n = \rho_n(X_k, Y) \sqrt{\frac{n-2}{1-\rho_n^2(X_k, Y)}}$  and  $\tilde{t}_n = \tilde{\rho}_n(X_k, Y) \sqrt{\frac{n-2}{1-\tilde{\rho}_n^2(X_k, Y)}}$  follow Student distributions with  $n-2$  degrees of liberty. From this, **significance tests can easily be proposed for screening in a sensitivity analysis context.**

Despite of their simple formulations, the Pearson's and Spearman's coefficients take into account only linear and monotonous effects respectively. Consequently, **they cannot deal with non-monotonic behavior and interactions between input parameters.**

## 2.2 The distance correlation

To address the limitations of correlation coefficients listed in Section 2.1, a first dependence measure presented in Da Veiga (2014) offers an interesting alternative. This quantity is based on the marginal distributions of the couple  $(X_k, Y)$  and avoids making parametric assumptions on the model  $Y = f(X)$ . Considering the random variables  $X_k \in \mathcal{X}_k \subset \mathbb{R}^{d_k}$  and  $Y \in \mathcal{Y} \subset \mathbb{R}^p$  with characteristic functions  $\Phi_{X_k}$  and  $\Phi_Y$ , the distance covariance is defined by

$$\mathcal{V}^2(X_k, Y) = \int_{\mathbb{R}^{d_k+p}} |\Phi_{X_k, Y}(t, s) - \Phi_{X_k}(t) \Phi_Y(s)|^2 w(t, s) dt ds \quad (1)$$

where  $w(t, s) = (c_{d_k} c_p \|t\|_2^{1+d_k} \|s\|_2^{1+p})^{-1}$  with the constant  $c_l = \pi^{(1+l)/2} / \Gamma((1+l)/2)$  for  $l \in \mathbb{N}$  and  $\|\cdot\|_2$  is the  $L^2$  norm (Székely et al, 2007). This quantity  $\mathcal{V}^2(X_k, Y)$  is equal to zero if and only if the characteristic function  $\Phi_{X_k, Y}$  of  $(X_k, Y)$  is equal to the product of  $\Phi_{X_k}$  and  $\Phi_Y$ , that is to say only and only if  $X_k$  and  $Y$  are independent. In other words, **the distance covariance is a good indicator of the dependence between  $X_k$  and  $Y$ , without any hypothesis on the law of  $X_k$  or the type of relation between  $X_k$  and  $Y$ .**

This distance covariance (1) can be expressed in terms of Euclidean distances:

$$\begin{aligned} \mathcal{V}^2(X_k, Y) &= \mathbb{E}_{X_k, X'_k, Y, Y'} [\|X_k - X'_k\|_2 \|Y - Y'\|_2] \\ &+ \mathbb{E}_{X_k, X'_k} [\|X_k - X'_k\|_2] \mathbb{E}_{Y, Y'} [\|Y - Y'\|_2] \\ &- 2 \mathbb{E}_{X_k, Y} [\mathbb{E}_{X'_k} [\|X_k - X'_k\|_2] \mathbb{E}_{Y'} [\|Y - Y'\|_2]] \end{aligned}$$

where  $(X', Y')$  is an independent and identically distributed copy of  $(X, Y)$  and where  $\mathbb{E}_Z$  represents the statistical mean in  $Z$ , for any random variable  $Z$ . From this statement, Székely et al (2007) propose an estimator of  $\mathcal{V}^2(X_k, Y)$ :

$$\begin{aligned} \mathcal{V}_n^2(X_k, Y) &= \frac{1}{n^2} \sum_{i,j=1}^n \|X_k^{(i)} - X_k^{(j)}\|_2 \|Y^{(i)} - Y^{(j)}\|_2 \\ &+ \frac{1}{n^4} \sum_{i,j=1}^n \|X_k^{(i)} - X_k^{(j)}\|_2 \sum_{i,j=1}^n \|Y^{(i)} - Y^{(j)}\|_2 \\ &- \frac{2}{n^3} \sum_{i=1}^n \left[ \sum_{j=1}^n \|X_k^{(i)} - X_k^{(j)}\|_2 \sum_{j=1}^n \|Y^{(i)} - Y^{(j)}\|_2 \right] \end{aligned}$$

This formulation can be condensed:  $\mathcal{V}_n^2(X_k, Y) = \frac{1}{n^2} \text{Tr} [G^{\{k\}} H G H]$  where  $H = (\delta_{ij} - \frac{1}{n})_{1 \leq i, j \leq n}$ ,  $G_k^{\{k\}} = (\|X_k^{(i)} - X_k^{(j)}\|_2)_{1 \leq i, j \leq n}$  and  $G = (\|Y^{(i)} - Y^{(j)}\|_2)_{1 \leq i, j \leq n}$ .

Székely et al (2007) propose another writing of  $\mathcal{V}_n^2(X_k, Y)$ , which is computationally cheaper:  $\mathcal{V}_n^2(X_k, Y) = \frac{1}{n^2} \sum_{i,j=1}^n A_{ij} B_{ij}$ , where  $A_{ij} = G_{ij}^{\{k\}} - \bar{G}_{i.}^{\{k\}} - \bar{G}_{.j}^{\{k\}} + \bar{G}_{..}^{\{k\}}$  and  $B_{ij} = G_{ij} - \bar{G}_{i.} - \bar{G}_{.j} + \bar{G}_{..}$  with  $\bar{M}_{.j} = n^{-1} \sum_{i=1}^n M_{ij}$ ,  $\bar{M}_{i.} = n^{-1} \sum_{j=1}^n M_{ij}$  and  $\bar{M}_{..} = n^{-2} \sum_{i,j=1}^n M_{ij}$ , for all  $M \in \mathcal{M}_n(\mathbb{R})$ .

Finally from the distance covariance, a distance correlation  $\mathcal{R}^2(X_k, Y)$  is proposed by Da Veiga Da Veiga (2014), defined by  $\mathcal{R}^2(X_k, Y) = \frac{\mathcal{V}^2(X_k, Y)}{\sqrt{\mathcal{V}^2(X_k, X_k) \mathcal{V}^2(Y, Y)}}$  if  $\mathcal{V}^2(X_k, X_k) \mathcal{V}^2(Y, Y) > 0$  and 0 otherwise. **This sensitivity index  $\mathcal{R}^2$  is included in the interval  $[0, 1]$ , like the absolute Pearson's correlation coefficient, which makes its interpretation easier.** The associated plug-in estimator deduced from  $\mathcal{V}_n^2(X_k, Y)$  is  $\mathcal{R}_n^2(X_k, Y) = \frac{\mathcal{V}_n^2(X_k, Y)}{\sqrt{\mathcal{V}_n^2(X_k, X_k) \mathcal{V}_n^2(Y, Y)}}$ .

## 2.3 The Hilbert-Schmidt dependence measure

Instead of quantifying the link between an input parameter and the model output from an analysis of their characteristic functions, Gretton et al (2005) propose to use the covariance between some transformations of these random variables. More precisely, they consider the random variables  $X \in \mathcal{X}_k$  and  $Y \in \mathcal{Y}$ , with the probability density functions  $p_{X_k}$  and  $p_Y$  and where  $\mathcal{X}_k$  and  $\mathcal{Y}$  are any measurable spaces. Then, they associate to  $X_k$  an universal Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{F}_k$  composed of functions mapping from  $\mathcal{X}_k$  to  $\mathbb{R}$  and defined by the kernel function  $\kappa_k$  (Aronszajn, 1950). The same transformation is realized with  $Y$ , considering the universal RKHS  $\mathcal{G}$  and the kernel function  $\kappa$ .  $\langle \cdot, \cdot \rangle_{\mathcal{F}_k}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{G}}$  are the scalar product over  $\mathcal{F}_k$  and  $\mathcal{G}$  respectively.

Then, the operator of crossed-covariance  $C_{X_k Y}$  associated to the probability density function  $p_{X_k Y}$  of  $(X_k, Y)$  is the linear operator mapping from  $\mathcal{G}$  to  $\mathcal{F}_k$  and defined for all  $f \in \mathcal{F}_k$  and for all  $g \in \mathcal{G}$  by  $\langle f, C_{X_k Y} g \rangle_{\mathcal{F}_k} = \text{Cov}(f(X_k), g(Y))$ . **This operator generalizes the covariance matrix between  $X_k$  and  $Y$ .** Indeed, thanks to the non-linear kernels which remove hypotheses such as linearity or monotony, it takes into account dependences more complex than the Pearson's and Spearman's coefficients.

Finally, the Hilbert-Schmidt Independence Criterion (HSIC) is the Hilbert-Schmidt norm of the operator  $C_{X_k Y}$  (Deza and Deza, 2009) :  $\text{HSIC}(X_k, Y) = \|C_{X_k Y}\|_{HS}^2 = \sum_{i,j} \langle u_i, C_{X_k Y} v_j \rangle_{\mathcal{F}_k}$  where  $(u_i)_{i \geq 0}$  and  $(v_j)_{j \geq 0}$  are orthonormal bases of  $\mathcal{F}_k$  and  $\mathcal{G}$ , respectively (Gretton et al, 2005). More precisely, the HSIC is equal to:

$$\begin{aligned} \text{HSIC}(X_k, Y) &= \mathbb{E}_{X_k, X'_k, Y, Y'} [\kappa_k(X_k, X'_k) \kappa(Y, Y')] \\ &+ \mathbb{E}_{X_k, X'_k} [\kappa_k(X_k, X'_k)] \mathbb{E}_{Y, Y'} [\kappa(Y, Y')] \\ &- 2 \mathbb{E}_{X_k, Y} [\mathbb{E}_{X'_k} [\kappa_k(X_k, X'_k)] \mathbb{E}_{Y'} [\kappa(Y, Y')]]. \end{aligned} \quad (2)$$

**Similarly to the distance covariance, the HSIC is equal to zero if and only if  $X_k$  and  $Y$  are independent**, without emitting any hypothesis about the nature of the relation between  $X_k$  and  $Y$ . In GSA, this property can be useful for a screening purpose. Moreover, for quantitative GSA, the HSIC can be used directly for quantifying uncertainty sources or normalized; Da Veiga (2014) proposes a sensitivity index in this direction:  $\mathcal{R}^2 \text{HSIC}(X_k, Y) = \frac{\text{HSIC}(X_k, Y)}{\sqrt{\text{HSIC}(X_k, X_k) \text{HSIC}(Y, Y)}}$ .

From a  $n$ -sample  $(X^{(i)}, Y^{(i)})_{1 \leq i \leq n}$  of  $(X, Y)$ , an estimator of the measure  $\text{HSIC}(X_k, Y)_{\mathcal{F}_k, \mathcal{G}}$  is proposed in (Gretton et al, 2005) :  $\text{HSIC}_n(X_k, Y) = \frac{1}{n^2} \text{Tr}(K_k H K H)$ , where  $K_k = (\kappa_k(X_k^{(i)}, X_k^{(j)}))_{1 \leq i, j \leq n}$  and  $K = (\kappa(Y^{(i)}, Y^{(j)}))_{1 \leq i, j \leq n}$ . Following the same way as Székely et al (2007), we propose to reduce the estimator calculation time using the previous formulation  $\frac{1}{n^2} \sum_{i,j=1}^n A_{ij} B_{ij}$  where now,  $G^{\{k\}} := K_k$  and  $G := K$ .

The kernel functions involved in the HSIC definition can belong to various classes of kernel functions, such as the Gaussian, the Laplacian or the Matérn family (Fukumizu et al, 2009). Note that these functions often require hyperparameter values which can be deduced from heuristic processes or fixed in order to maximize the HSIC value (Balasubramanian et al, 2013). In this paper, we consider the Gaussian kernel function  $k(z^{(i)}, z^{(j)}) = \exp\left(-\sum_{k=1}^{n_z} \frac{(z_k^{(i)} - z_k^{(j)})^2}{\sigma_k^2}\right)$  for inputs and outputs and  $\sigma^2$  is estimated by the empirical variance associated to  $z_k^{(1)}, \dots, z_k^{(n)}$  (Yamada et al, 2014).

### 3 Significance tests for screening purpose

In a screening context, the objective is to separate the input parameters into two sub-groups, the significant ones and the non-significant ones. For this, we propose to use statistical hypothesis tests based on dependence measures described in Section 2. For a given input  $X_k$ , it aims at testing the null hypothesis " $\mathcal{H}_0^{(k)}$ :  $X_k$  and  $Y$  are independent", against its alternative " $\mathcal{H}_1^{(k)}$ :  $X_k$  and  $Y$  are dependent". The significance level<sup>1</sup> of these tests is hereinafter noted  $\alpha$ . Some asymptotic results exist in this domain for the dependence measures; we briefly present some of them in the following, based on the notations of Section 2. In a second part, we develop spectral approximations of the asymptotic laws governing the statistics involved in these tests, which can be useful for medium size samples. Finally, we propose to extend these results to the non-asymptotic case thanks to a bootstrap approach.

#### 3.1 Asymptotic tests of independence

##### Asymptotic test for the HSIC

Considering the HSIC, Gretton et al (2007) propose a kernel statistical test of independence based on asymptotic considerations. First, the estimator  $\text{HSIC}_n(X_k, Y)$  is rewritten:  $\text{HSIC}_n(X_k, Y) = \frac{1}{n^4} \sum_{i,j,q,r} h_{ijqr}$  where  $h_{ijqr} = \frac{1}{4!} \sum_{(t,u,v,w)}^{(i,j,q,r)} K_{k,tu} K_{tu} + K_{k,tu} K_{vw} - 2K_{k,tu} K_{tv}$ , the sum being done over the different permutations  $(t, u, v, w)$  of  $(i, j, q, r)$ . Then, under  $\mathcal{H}_0$ , the statistic  $n\text{HSIC}_n(X_k, Y)$  converges in distribution to  $\sum_{l>0} \lambda_l Z_l^2$ , where the standard normal variables  $Z_l$  are independent and where the coefficients  $\lambda_l$  are the solutions of the eigenvalues problem  $\lambda_l \psi_l(z_j) = \int h_{ijqr} \psi_l(z_i) dF_{iqr}$ ,  $F_{iqr}$  being the distribution function of  $(Z_i, Z_q, Z_r)$  and  $\psi_l(\cdot)$  the eigenvector associated to  $\lambda_l$ .

In practice (for details, see Gretton et al, 2007), the distribution of the infinite weighted sum of independent chi-squared variables is approached by a Gamma distribution with shape parameter  $\gamma$  and inverse scale parameter  $\beta$ . These parameters are estimated by  $\hat{\gamma} = \frac{n^{-2}(1+E_x E_y - E_x - E_y)^2}{V}$  and  $\hat{\beta} = \frac{nV}{n^{-1}(1+E_x E_y - E_x - E_y)}$  where  $E_x = \frac{1}{n(n-1)} \sum_{i \leq j \leq n} (K_k)_{ij}$ ,  $E_y = \frac{1}{n(n-1)} \sum_{i \leq j \leq n} K_{ij}$  and  $V = \frac{2(n-4)(n-5)}{n(n-1)(n-2)(n-3)} \mathbf{1}^T (B - \text{diag}(B)) \mathbf{1}$ , with  $B = ((HK_k H) \odot (HKH))^2$ .  $\odot$  is the element-wise multiplication and  $M^{\cdot 2}$  the element-wise matrix power for all  $M \in \mathcal{M}_n(\mathbb{R})$ .

Finally, the independence test rejects the null hypothesis  $\mathcal{H}_0$  when the  $p$ -value of the Gamma distribution associated to the statistic  $n\text{HSIC}_n(X_k, Y)$  is lower than some level  $\alpha$ , e.g.  $\alpha = 5\%$ .

##### Asymptotic test for the distance covariance

For the distance covariance introduced in Section 2.2, we refer to Székely et al (2007) in the case where  $\mathbb{E}[\|X_k\|_{d_k} + \|Y\|_p] < \infty$ .

Firstly, if  $X_k$  and  $Y$  are independent,  $\frac{nV_n^2}{S_2} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \sum_{l>0} \lambda_l Z_l^2$ , where the standard normal variables  $(Z_l)_{l>0}$  are independent and the  $\lambda_l$  are positive reals, with  $S_2 = n^{-2} \left( \sum_{i,j=1}^n G_{ij}^{(k)} \right) \left( \sum_{i,j=1}^n G_{ij} \right)$ . If  $X_k$  and  $Y$  are dependent,  $nV_n^2/S_2 \xrightarrow[n \rightarrow \infty]{\mathcal{P}} \infty$ .

Secondly, we consider  $T(X_k, Y, \alpha, n)$  the statistical test rejecting the null hypothesis " $\mathcal{H}_0$ :  $X_k$  and  $Y$  are independent" when  $\frac{nV_n^2}{S_2} > (\Phi^{-1}(1 - \alpha/2))^2$ ,  $\Phi$  being the distribution function of the standard normal law. If  $\mathbb{E}[\|X_k\|_{d_k} + \|Y\|_p] < \infty$ , then for all  $\alpha \in ]0, 0.215]$ ,  $\lim_{n \rightarrow \infty} \alpha(X_k, Y, n) \leq \alpha$  and  $\sup_{X_k, Y} \{\lim_{n \rightarrow \infty} \alpha(X_k, Y, n) : \mathcal{V}(X_k, Y) = 0\} = \alpha$ , where  $\alpha(X_k, Y, n)$  is the type I error rate of  $T(X_k, Y, \alpha, n)$ .

Consequently, the test  $T(X_k, Y, \alpha, n)$  has an asymptotic type I error rate at worst equal to  $\alpha$  and the approximation of the  $1 - \alpha$  quantile of the law of  $\sum_{l>0} \lambda_l Z_l^2$  by the squared  $1 - \alpha/2$  quantile of the standard normal law seems to be a powerful technique.

#### 3.2 Spectral approach for the asymptotic tests

For small and medium size samples, the previous approximations of the asymptotic laws are questionable. For example, Székely et al (2007) show that in the case of the distance covariance, the test  $T(X_k, Y, \alpha, n)$  might be over-conservative. In the context of a two-sample test, Shen et al (2009) remind us of the heuristic nature of the Gamma approximation for the asymptotic law of the HSIC estimator. This substitution of laws can be not enough accurate for the upper tail of the distribution, that is to say for its most important part in the case of a  $p$ -value computation. Consequently, Sejdinovic et al (2013) advise the use of a spectral approximation of the asymptotic

<sup>1</sup>The significance level of a statistical hypothesis test is the rate of the type I error which corresponds to the rejection of the null hypothesis  $\mathcal{H}_0$  when it is true.

laws for the HSIC and the distance covariance, which are weighted sums of chi-squares as mentioned in Section 3.1. Particularly, Zhang et al (2011) propose a spectral estimation of the asymptotic law of the HSIC estimator. Precisely, the asymptotic law of  $\frac{\text{HSIC}_n(X_k, Y)}{n}$  can be approached by those of  $\frac{1}{n^2} \sum_{i,j=1}^n \hat{\lambda}_{k,i} \hat{\nu}_j \varepsilon_{ij}^2$ , where  $\varepsilon_{ij}$ ,  $1 \leq i, j \leq n$  are independent standard normal variables and  $(\hat{\lambda}_{k,i})_{1 \leq i \leq n}$  and  $(\hat{\nu}_i)_{1 \leq i \leq n}$  are the eigenvalues of  $HG^{\{k\}}H$  and  $HGH$  respectively. This result can easily be extended to the distance covariance, replacing the statistic  $\frac{\text{HSIC}_n(X_k, Y)}{n}$  by  $nV_n^2/S_2$ .

As it requires only the computation of the matrix-vector product  $\hat{\lambda}' \varepsilon_n \hat{\nu}$  where  $\hat{\lambda} = (\hat{\lambda}_1, \dots, \hat{\lambda}_n)'$ ,  $\hat{\nu} = (\hat{\nu}_1, \dots, \hat{\nu}_n)'$  and  $\varepsilon_n = (\varepsilon_{ij})_{1 \leq i, j \leq n}$ , an instance of such random variables is clearly cheaper than a bootstrapped instance of the corresponding dependence measures. However this last approach can be required for small samples.

### 3.3 Non-asymptotic tests based on resampling

The significance tests based on dependence measures presented in Sections 3.1 and 3.2 are fast and asymptotically efficient tools for the selection of the influential input parameters. However, they are considerably biased when the number of observations  $n$  is too weak because of their asymptotic framework. Consequently, non-asymptotic results are necessary.

For this purpose, we propose a generic non-parametric test based on resampling, which can be applied to any dependence measure  $\Delta(X_k, Y)$  between two random variables  $X_k$  and  $Y$ . For this,  $B$  bootstrap versions  $\mathbf{Y}^{[1]}, \dots, \mathbf{Y}^{[B]}$  of the output sample  $\mathbf{Y} = (Y^{(1)} \dots Y^{(n)})$  are generated and for each  $\mathbf{Y}^{[b]}$ , the associated input sample is  $\mathbf{X}_k^{[b]} := \mathbf{X}_k$  where  $\mathbf{X}_k = (X_k^{(1)} \dots X_k^{(n)})$ .

Under these considerations, our test can be summarized by the following algorithm:

1. Create a sample  $(\mathbf{X}, \mathbf{Y}) = (X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)})_{1 \leq i \leq n}$ .
2. Compute  $\Delta_n(X_k, Y)$ , an estimator of the dependence measure  $\Delta(X_k, Y)$ .
3. Realize  $B$  bootstrap samplings  $(\mathbf{X}_k^{[b]}, \mathbf{Y}^{[b]})$  of the sample  $(\mathbf{X}_k, \mathbf{Y})$  under  $\mathcal{H}_0$ .
4. Compute  $(\Delta_n(X_k^{[b]}, Y^{[b]}))_{1 \leq b \leq B}$ , the  $B$  bootstrap estimations.
5. Compute the bootstrapped  $p$ -value  $p\text{-val}_B = \frac{1}{B} \sum_{b=1}^B \mathbb{1}_{\Delta_n(X_k^{[b]}, Y^{[b]}) > \Delta_n(X_k, Y)}$ .
6. If  $p\text{-val}_B < \alpha$ , then reject  $\mathcal{H}_0$ , else accept  $\mathcal{H}_0$ .

**Remark 1.** This algorithm is designed for testing the dependence between an input parameter and the output of the model. If we want to simultaneously apply this test for the  $d$  input parameters, only the steps 4 to 6 have to be repeated for each input. This avoids  $d - 1$  repetitions of the bootstrap step.

### 3.4 Synthesis on significance tests

For methodological recommendations:

- We propose to use the significance tests based on resampling (Section 3.3) in presence of a small sample.
- When the number of observations is much more important, we advise to use the asymptotic tests (Section 3.1).
- Between both situations, we propose to use the spectral approach (Section 3.2), which is better than both the approximation of the asymptotic laws and the use of the empirical distribution of a dependence measure estimator. Indeed, even if this last law is more justified than the asymptotic one, Sejdinovic et al (2013) highlight its important cost. This is due to the computation of the dependence measure estimator for each bootstrapped sample, especially when the input or output parameter space dimension is important.

## 4 Bootstrapped linear regression for dependence measures in screening

The previous significance tests for feature selection are directly computed on the dependence measures presented in Section 2, which associate one input parameter to the model output. In this section, we propose to decompose in a linear way the difference between two output observations according to the differences between the associated input observations; we call “local measures” these simple quantities measuring the difference between two observations of a same variable. Considering this linear model, our aim is to build significance tests for the different effects, using classical tools for nested model selection. Discarding an effect from this regression model corresponds to discarding a significant input parameter in a screening context.

### 4.1 Linear model between the local measures

Considering a  $n$ -sample  $(X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)})_{1 \leq i \leq n}$  and a local measure  $D(\cdot, \cdot)$ , we propose the linear model:

$$D(Y^{(i)}, Y^{(j)}) = \sum_{k=1}^d \beta_k D(X_k^{(i)}, X_k^{(j)}), \quad 1 \leq i, j \leq n \quad (3)$$

where  $\beta \in \mathbb{R}_+^d$ . For two observations  $(X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)})$  and  $(X_1^{(j)}, \dots, X_d^{(j)}, Y^{(j)})$ , the coefficient  $\beta_k$  can be interpreted as the weight associated to the contribution of the dependence between  $X_k^{(i)}$  and  $X_k^{(j)}$  to the explanation of the dependence between  $Y^{(i)}$  and  $Y^{(j)}$ .

The vector  $\beta$  can be estimated by:

$$\hat{\beta} \in \underset{\beta \in (\mathbb{R}_+)^d}{\operatorname{arginf}} \left\| D(\mathbf{Y}) - \sum_{k=1}^d \beta_k D(\mathbf{X}_k) \right\|_{\text{Frob}}^2 \quad (4)$$

where  $\|\cdot\|_{\text{Frob}}$  is the Frobenius norm defined for all  $A \in \mathcal{M}_n(\mathbb{R})$  by  $\|A\|_{\text{Frob}} = \sqrt{\sum_{i,j=1}^n A_{ij}^2}$  and where  $D(\mathbf{A}) = (D(A_i, A_j))_{1 \leq i, j \leq n}$ . As a function of the random variables  $\mathbf{X}$  and  $\mathbf{Y}$ ,  $\hat{\beta}$  is also a random variable. For an easier implementation, we can rewrite Equation (4) with the Euclidean norm, replacing the matrix evaluations of the local measure  $D$  by their vectorized forms:

$$\hat{\beta} \in \underset{\beta \in (\mathbb{R}_+)^d}{\operatorname{arginf}} \left\| \vec{D}(\mathbf{Y}) - \sum_{k=1}^d \beta_k \vec{D}(\mathbf{X}_k) \right\|_2^2 = \underset{\beta \in (\mathbb{R}_+)^d}{\operatorname{arginf}} \left\| \vec{D}(\mathbf{Y}) - [\vec{D}(\mathbf{X}_1) \dots \vec{D}(\mathbf{X}_d)] \beta \right\|_2^2.$$

where  $(\vec{D}(\mathbf{Y}))_{(j-1)n+i} := D(Y^{(i)}, Y^{(j)}), \forall i, j \in \{1, \dots, n\}$ , and so on.

**Remark 2.** In practice, the symmetric property of the matrices  $D(\mathbf{Y})$ ,  $D(\mathbf{X}_1)$ , ... and  $D(\mathbf{X}_d)$  allows the use of smaller vectors  $\vec{D}(\mathbf{Y})$ ,  $\vec{D}(\mathbf{X}_1)$ , ... and  $\vec{D}(\mathbf{X}_d)$  of size  $\frac{n(n+1)}{2}$  instead of  $n^2$ .

**Remark 3.** The decomposition of the  $Y$  local measure into a linear combination of  $X_1, \dots, X_d$  local measures makes sense if the coefficients  $\beta_1, \dots, \beta_d$  are non-negative. This is the reason why the problem is a constrained linear least-squares minimization with  $\beta \in (\mathbb{R}_+)^d$  rather than a simple linear least-squares minimization with  $\beta \in \mathbb{R}^d$ . This consideration leads to a more expensive problem resolution because of numerical optimization steps instead of an analytical solution  $\hat{\beta}$ .

The objective function in the constrained minimization problem (4) takes the form  $\eta(\mathbf{X}, \mathbf{Y}; \beta) = \left\| D(\mathbf{Y}) - \sum_{k=1}^d \beta_k D(\mathbf{X}_k) \right\|_{\text{Frob}}^2$  and can be decomposed in the sum of three terms:

$$\Delta(\mathbf{Y}, \mathbf{Y}) - 2 \sum_{k=1}^d \beta_k \Delta(\mathbf{X}_k, \mathbf{Y}) + \sum_{k,l=1}^d \beta_k \beta_l \Delta(\mathbf{X}_k, \mathbf{X}_l), \quad (5)$$

where  $\Delta(\mathbf{A}, \mathbf{B}) = \operatorname{Tr}[D(\mathbf{A})D(\mathbf{B})^T] \geq 0$ . For certain local measures  $D$ ,  $\Delta$  quantifies the global dependence between the random variables  $A$  and  $B$  using  $n$  independent evaluations stocked in  $\mathbf{A}$  and  $\mathbf{B}$ . In these particular cases, the proposed scheme (5) is linked to the “minimal-redundancy-maximal-relevance” strategy (mRMR) because its minimization gives important weights to the input parameters maximizing the dependence measures  $\Delta(\mathbf{X}_k, \mathbf{Y})$  and small weights to the input parameters highly dependent to the previous ones (Peng et al, 2005). This can be very useful when many input parameters are dependent: in the extreme case where a parameter input is no more than a deterministic function of another one, we would be interested in a method keeping only one of these two variables.



Especially, in the case of the HSIC and using the notations of Section 2.3, the choices  $D(\mathbf{Y}) = HKH$  and  $D(\mathbf{X}_k) = HK_kH$ ,  $\forall k \in \{1, \dots, d\}$ , lead to a result presented in Da Veiga (2014):  $\eta(\mathbf{X}, \mathbf{Y}; \beta) = \text{HSIC}_n(Y, Y) - 2 \sum_{k=1}^d \beta_k \text{HSIC}_n(X_k, Y) + \sum_{k,l=1}^d \beta_k \beta_l \text{HSIC}_n(X_k, X_l)$ . Likewise, in the case of the distance covariance and using the notations of Section (2.2), the choices  $D(\mathbf{Y}) = HGH$  and  $D(\mathbf{X}_k) = HG^{(k)}H$ ,  $\forall k \in \{1, \dots, d\}$ , lead to the mRMR scheme  $\eta(\mathbf{X}, \mathbf{Y}; \beta) = \mathcal{V}_n^2(Y, Y) - 2 \sum_{k=1}^d \beta_k \mathcal{V}_n^2(X_k, Y) + \sum_{k,l=1}^d \beta_k \beta_l \mathcal{V}_n^2(X_k, X_l)$ . In a similar way, for the Pearson's coefficient correlation  $\rho(X_k, Y)$ , we can show that the choices of  $D(\mathbf{Y}) = HYY^T H$  and  $D(\mathbf{X}_k) = HX_kX_k^T H$  lead to the mRMR scheme  $\eta(\mathbf{X}, \mathbf{Y}; \beta) = \text{Cov}_n^2(Y, Y) - 2 \sum_{k=1}^d \beta_k \text{Cov}_n^2(X_k, Y) + \sum_{k,l=1}^d \beta_k \beta_l \text{Cov}_n^2(X_k, X_l)$ .

In the following, we consider an alternative to the estimator (4) for the estimation of the regression parameters in the linear model (3), particularly useful when the number of input parameters is important.

## 4.2 Shrinkage in high-dimension

The coefficient estimation in the linear model (3) can be realized using regularization techniques. These methods are said active because they select the optimal complexity of the model during the optimization step (4) modified in some manner. More precisely, these techniques consist in the minimization of a quadratic risk penalized by an additive term, which is a constraint on the number or the size of model parameters, such as a limited  $\ell^2$  norm (Hoerl and Kennard, 1970),  $\ell^1$  norm (Tibshirani, 1996) or a combination of both (Zou and Hastie, 2005). In other words, a shrinkage tool looks for the optimal parameter values of (3) and the optimal effective dimension of the problem.

Under these considerations, we could use the Lasso (Least Absolute Shrinkage and Selection Operator) penalty (Tibshirani, 1996) in order to select a subset of the local measures in the full model (3), and so a subset of the input parameters:

$$\eta_{\text{lasso}}(\mathbf{X}, \mathbf{Y}; \beta) = \left\| D(\mathbf{Y}) - \sum_{k=1}^d \beta_k D(\mathbf{X}_k) \right\|_{\text{Frob}}^2 + \lambda \|\beta\|_1.$$

It is in this sense that Yamada et al (2014) propose the HSIC Lasso which consists in the minimization of this objective function with  $D(\mathbf{Y}) = HKH$  and  $D(\mathbf{X}_k) = HK_kH$ ,  $\forall k \in \{1, \dots, d\}$ , under the positivity constraints  $\beta_1 \geq 0, \dots, \beta_d \geq 0$  and using a dual augmented Lagrangian algorithm to solve the optimization problem.

In this paper, the HSIC Lasso is used but, for time computation reasons, we propose to solve the optimization problem with the Least Angle Regression (LARS) algorithm under positivity constraints (Efron et al, 2004, Sec. 3.4), with a regularization parameter  $\lambda$  optimized by an improved version of the cross-validation error minimization. Usually we take  $\hat{\lambda}_{\text{CV}}$ , the  $\lambda$  value minimizing the cross-validation error  $\mu_{\text{CV}}^{(l)}$ . In the HSIC lasso framework, we propose to replace  $\hat{\lambda}_{\text{CV}}$  by  $\hat{\lambda}_{\text{CV mod}}$  which minimizes

$$\mu_{\text{CV}}^{(l)} - 0.5\sigma_{\text{CV}}^{(l)}$$

over the indices  $\{1, \dots, L\}$  of the discretized  $\lambda$  values for the optimization,  $\sigma_{\text{CV}}^{(l)}$  being the standard deviation of the prediction error associated to the different folds. The 0.5 value has been chosen after tests over various analytical functions.  $\mu_{\text{CV}}^{(l)} - 0.5\sigma_{\text{CV}}^{(l)}$  is an amelioration of the  $\mu_{\text{CV}}^{(l)}$  minimization because it takes into account the uncertainty of the cross-validation error.

Finally, if the effective dimension of the problem is of the same order than the number of input parameters, nested model selection tools can be considered instead of the shrinkage approach.

## 4.3 Bootstrap test for the nested model selection

Based on the full model (3) of Section 4.1, we propose some methods using significance tests in order to remove the non-significant input parameters. More precisely, for a given input parameter  $X_k$ , we want to build a statistical test with the null hypothesis " $\mathcal{H}_0^{(k)}$ :  $X_k$  and  $Y$  are independent" and its alternative " $\mathcal{H}_1^{(k)}$ :  $X_k$  and  $Y$  are dependent", or in an equivalent way: " $\mathcal{H}_0^{(k)}$ :  $\beta_k = 0$ " and " $\mathcal{H}_1^{(k)}$ :  $\beta_k \neq 0$ ". Obviously, the law of  $\hat{\beta}_k$  in (4) is unknown and, at best, asymptotically approximable. Consequently, similarly to the resampling method proposed in Section 3.3, we propose to build a bootstrap test for each input parameter  $X_k$ , starting from the  $n$ -sample  $(X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)})_{1 \leq i \leq n}$ .

More precisely, the  $b^{\text{th}}$  bootstrap sample  $(X_1^{[b],(i)}, \dots, X_d^{[b],(i)}, Y^{[b],(i)})_{1 \leq i \leq n}$  is such that

$$Y^{[b],(i)} := Y^{(i)}, X_l^{[b],(i)} := X_l^{(i)}, \forall l \neq k \text{ and } X_k^{[b],(i)} := X_k^{[b],(i)}.$$

In other words, the  $b^{\text{th}}$  bootstrap sample corresponds to the  $n$ -sample  $(X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)})_{1 \leq i \leq n}$  where the observations of the  $k^{\text{th}}$  input parameter are resampled according to their empirical probability distribution. Then, we

compute  $\hat{\beta}^{[b]}$ , the estimation of the vector  $\beta$  for each bootstrap sample  $(X_1^{[b],(i)}, \dots, X_d^{[b],(i)}, Y^{[b],(i)})_{1 \leq i \leq n}$ .

Afterwards, under the null hypothesis  $\mathcal{H}_0^{(k)}$ , the  $p$ -value is estimated by  $p\text{-val}_B^{(k)} = \frac{1}{B} \sum_{b=1}^B \mathbb{1}_{\hat{\beta}_k^{[b]} > \hat{\beta}_k}$  and  $\mathcal{H}_0^{(k)}$  is rejected when  $p\text{-val}_B^{(k)}$  is lower than some level  $\alpha$ , e.g.  $\alpha = 5\%$ .

Finally, considering the conclusions of the  $d$  statistical tests, we obtain a sub-model of the full model (4) keeping only the significant local measures, and from another point of view, dismissing the non-significant input parameters of the model  $Y = f(X_1, \dots, X_d)$ . From this conclusion, we could imagine to go further and to apply the tools commonly used for feature selection in the linear model, such as forward, backward or stepwise approaches.

## 5 Numerical experiments

In these sections, we numerically investigate the methods expounded in Sections 3 and 4 for a screening purpose. We also compare the distance correlation and the HSIC with the classical Sobol' indices, in order to exhibit some of their specificities.

**Reminder on Sobol' indices.** For a model  $Y(X) = f(X_1, \dots, X_d) \in \mathbb{R}$  with independent random real variables  $X_1, \dots, X_d$  and such that  $\mathbb{E}[f^2(X)] < +\infty$ , we can apply the Hoeffding decomposition:

$$f(X) = f_0 + \sum_{i=1}^d f_j(X_j) + \sum_{i=1}^d \sum_{i < j}^d f_{ij}(X_i, X_j) + \dots + f_{1\dots d}(X_1, \dots, X_d) = \sum_{u \subset \{1, \dots, d\}} f_u(X_u)$$

where  $f_0 = \mathbb{E}[f(X)]$ ,  $f_j(X_j) = \mathbb{E}[f(X)|X_j] - f_0$  and  $f_u(X_u) = \mathbb{E}[f(X)|X_u] - \sum_{v \subset u} f_v(X_v)$ , with for all  $u \subset \{1, \dots, d\}$ ,  $X_u = (X_i)_{i \in u}$ . Then for each  $u \subset \{1, \dots, d\}$ , the first-order and total Sobol' indices of  $X_u$  are defined by  $S_u = \frac{\mathbb{V}[f_u(X_u)]}{\mathbb{V}[f(X)]}$  and  $S_u^T = \sum_{v \supset u} S_v$ , where  $\mu_{X_u}$  and  $\mu_X$  are the distribution functions of  $X_u$  and  $X$  respectively.

The first-order indices associated to  $X_1, \dots, X_d$  can also be rewritten:  $S_k = \frac{\mathbb{V}[\mathbb{E}[f(X)|X_k]]}{\mathbb{V}[f(X)]}$ ,  $\forall k \subset \{1, \dots, d\}$ .

### 5.1 Comparison of sensitivity indices

These first tests on analytical functions aim at comparing various sensitivity indices: the classical Sobol' indices vs. dependence measures such as distance correlation (dCor), HSIC and sup-HSIC (the supremum of HSIC over the possible correlation length values). **The objective is to identify which kinds of input effect they allow to detect, and to highlight any difference between these indices.**

For this, several analytical functions including linear or not, monotonic or not input effects are used in the following numerical tests. To build the different test functions, we considered monodimensional functions designed to be centered and with variance one when  $x$  is a realization of a uniform random variable on  $[-\sqrt{3}, \sqrt{3}]$ , centered with variance one. These elementary functions are of three type:

1. linear:  $h_1(x) = x$ ;
2. monotonous (exponential):  $h_2(x) = \frac{e^x - a}{b}$  where  $a = \frac{\sinh(\sqrt{3})}{\sqrt{3}}$  and  $b = \sqrt{\frac{\sinh(2\sqrt{3})}{2\sqrt{3}} - a^2}$ ;
3. non-monotonous (sinusoidal):  $h_3(x) = a \sin(2x)$  where  $a = 1/\sqrt{0.5 - \frac{\sin(4\sqrt{3})}{8\sqrt{3}}}$ .

The  $d$  input parameters  $X = (X_1, \dots, X_d)$  of model  $f$  are supposed independent and identically distributed according to a uniform distribution over  $[-\sqrt{3}, \sqrt{3}]$ .

#### Sensitivity indices regarding the shape of monodimensional effects

First of all, we consider the additive model  $f(X)$  with only monodimensional effects:  $f(X) = c_1 h_1(X_1) + c_2 h_2(X_2) + c_3 h_3(X_3)$ ,  $c \in \mathbb{R}^3$ , and we propose to study the sensitivity of Sobol, HSIC, sup-HSIC and dCor indices to linear, monotonous and non-monotonous effects. Note that, in this case, the total Sobol' indices are equal to the first-order ones:  $S_i = S_i^T = \frac{c_i^2}{c_1^2 + c_2^2 + c_3^2}$ ,  $i \in \{1, 2, 3\}$ . In the following tests, the coefficients  $c_i$  are set to 0 or 1, allowing to cancel the effect of the corresponding  $X_i$ .

The various sensitivity indices HSIC, sup-HSIC and dCor are estimated with a Monte-Carlo sampling of 1000 simulations, and compared to analytical Sobol index values. Note that, for this size of sampling, several Monte-Carlo repetitions have been performed and a negligible variance of dependence measure estimation has been observed, justifying this choice of sample size. For each index, estimation is repeated 100 times. The mean values of sensitivity indices obtained for each kind of model are given in Table 1 in percentage (the sensitivity index for an input parameter

$Y = f(X) = \dots$	Effect	HSIC	sup-HSIC	dCor	Sobol
$h_1(X_1) + h_2(X_2)$	$X_1$ : linear	<b>62</b>	<b>61</b>	<b>57</b>	50
	$X_2$ : monotonous	38	39	43	50
$h_1(X_1) + h_3(X_3)$	$X_1$ : linear	<b>55</b>	<b>55</b>	<b>63</b>	50
	$X_3$ : non-monotonous	45	45	37	50
$h_2(X_2) + h_3(X_3)$	$X_2$ : non-linear	44	45	<b>56</b>	50
	$X_3$ : non-monotonous	<b>56</b>	<b>55</b>	44	50
$h_1(X_1) + h_2(X_2) + h_3(X_3)$	$X_1$ : linear	<b>38</b>	<b>38</b>	<b>41</b>	33
	$X_2$ : non-linear	31	31	35	33
	$X_3$ : non-monotonous	31	31	24	33

Table 1: Sensitivity indices in percentage for different test functions.

is normalized by the sum of the sensitivity indices of different inputs). Firstly, the dependence measures HSIC, sup-HSIC and dCor are different from the Sobol' indices, with a relative difference up to 20% with respect to the latter. Secondly, HSIC and sup-HSIC give the same results for these test functions. Then, the dependence measures give additional weight to linear effects, in comparison with the Sobol' indices. Finally, we observe differences between HSIC (and sup-HSIC) and dCor for non-linear functions, HSIC highlighting non-monotonous effects while dCor featuring monotonous ones.

### Sensitivity indices regarding the weight of the interaction effect

Now, we consider the additive model  $f(X) = h_2(X_1) + ch_2(X_1)h_2(X_2)$  with a monodimensional and an interaction effect, the latter being weighted by a positive real  $c$ . We propose to study the sensitivity of Sobol and HSIC to the value of  $c$ . For brevity, we only present results for exponential shape, the conclusion being qualitatively the same for linear and sinusoidal ones. For the same reason, we only consider the dependence measure HSIC. The first-order Sobol' indices are  $S_1 = \frac{1}{1+c^2}$  and  $S_2 = 0$  while the total ones are equal to  $S_1^T = 1$  and  $S_2^T = \frac{c^2}{1+c^2}$  for all  $c \in \mathbb{R}_+$ .

Measure	$c$									
	0		1		2		4		5	
HSIC	$X_1$	$X_2$	$X_1$	$X_2$	$X_1$	$X_2$	$X_1$	$X_2$	$X_1$	$X_2$
HSIC <sub>k</sub> / $\sum_{j=1}^d \text{HSIC}_j$	<b>0.0965</b>	0.0003	0.0293	<b>0.0309</b>	0.0071	<b>0.0250</b>	0.0092	<b>0.0184</b>	0.0104	<b>0.0176</b>
Borgonovo $\delta_k$	<b>0.7759</b>	0.0044	0.4110	<b>0.4530</b>	0.2845	<b>0.3993</b>	0.2971	<b>0.3610</b>	0.3022	<b>0.3546</b>
$\delta_k / \sum_{j=1}^d \delta_j$	<b>99.5%</b>	0.6%	47.8%	<b>52.4%</b>	41.6%	<b>58.4%</b>	45.1%	<b>54.9%</b>	46.0%	<b>54.0%</b>
Total Sobol $S_k^T$	<b>1</b>	0	<b>1</b>	0.5000	<b>1</b>	0.8000	<b>1</b>	0.9412	<b>1</b>	0.9615
$S_k^T / \sum_{j=1}^d S_j^T$	<b>100%</b>	0%	<b>66.7%</b>	33.3%	<b>55.6%</b>	44%	<b>51.5%</b>	48.5%	<b>51.0%</b>	49.0%

Measure	$c$									
	6		7		8		9		10	
HSIC	$X_1$	$X_2$	$X_1$	$X_2$	$X_1$	$X_2$	$X_1$	$X_2$	$X_1$	$X_2$
HSIC <sub>k</sub> / $\sum_{j=1}^d \text{HSIC}_j$	0.0112	<b>0.0171</b>	0.0118	<b>0.0168</b>	0.0123	<b>0.0166</b>	0.0127	<b>0.0165</b>	0.0130	<b>0.0164</b>
Borgonovo $\delta_k$	39.6%	<b>60.4%</b>	41.3%	<b>58.7%</b>	42.3%	<b>57.4%</b>	43.5%	<b>56.5%</b>	44.2%	<b>55.8%</b>
$\delta_k / \sum_{j=1}^d \delta_j$	0.3059	<b>0.3497</b>	0.3086	<b>0.3460</b>	0.3106	<b>0.3432</b>	0.3121	<b>0.3410</b>	0.3133	<b>0.3392</b>
Total Sobol $S_k^T$	46.7%	<b>53.3%</b>	47.1%	<b>52.9%</b>	47.5%	<b>52.5%</b>	47.8%	<b>52.2%</b>	48.0%	<b>52.0%</b>
$S_k^T / \sum_{j=1}^d S_j^T$	<b>1</b>	0.9730	<b>1</b>	0.9800	<b>1</b>	0.9846	<b>1</b>	0.9878	<b>1</b>	0.9901
	<b>50.7%</b>	49.3%	<b>50.5%</b>	49.5%	<b>40.4%</b>	49.6%	<b>50.3%</b>	49.7%	<b>50.2%</b>	49.8%

Table 2: Standard and normalized HSIC, Borgonovo and total Sobol' indices for different values of  $c$ .

Table 2 presents the mean HSIC estimations associated to this study for different values of  $c$  based on 1000 repetitions of a 1000-sample, while Figure 1 illustrates the evolution of these indices according to the value of  $c$  for a certain 1000-sample. Naturally, the first variable is the more influential when  $c \ll 1$  because the second variable is almost missing in the model. Then, both variables tend to have the same effect around  $c = 1$  and finally the second variable is the more influential with a pick around  $c = 2$  where the HSICs associated to  $X_1$  and  $X_2$  start a convergence to the same value.

Table 2 and Figure 1 show that the same phenomenon occurs with the Borgonovo's delta moment independent measure (Plischke et al, 2013) defined by  $\delta_k = \frac{1}{2} \int_{\mathcal{X}_k} f_{X_k}(x) \int_{\mathcal{Y}} |f_Y(y) - f_{Y|X_k=x}(y)| dy dx$ , where  $f_{X_k}$ ,  $f_Y$  and  $f_{Y|X_k=x}$  are the density probability functions of  $X_k$ ,  $Y$  and  $Y|X_k$  respectively. Moreover, we have found the same result with the randomized dependence coefficient, a recent dependence measure defined in terms of correlation of random non-linear copula projections (López-Paz et al, 2013). In the framework of the classical output variance decomposition, this situation is surprising: the second variable occurring only in the interaction effect, its total

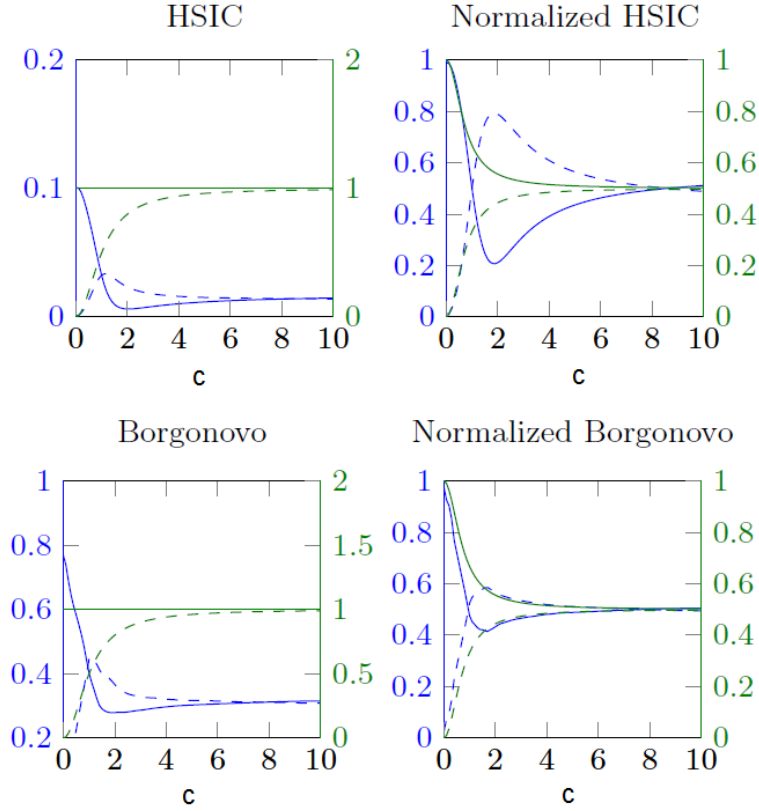


Figure 1: Standard and normalized HSIC and Borgonovo indices (blue lines) of  $X_1$  (plain lines) and  $X_2$  (dashed lines) for different values of  $c$ , with the corresponding total Sobol' indices (green lines).

Sobol' index  $S_2^T = c^2(1+c^2)^{-1}$  is obviously lower than the first input one  $S_1^T = 1$  for every  $c$ . However, such situation is understandable if we look at the model  $f$  from a multiplicative point of view, instead of an additional one:  $f(X) = h(X_1)(1+ch(X_2))$ . This leads to the corresponding multiplicative decomposition of the variance:  $\mathbb{V}[f(X)] = \mathbb{E}[(h(X_1))^2] \times \mathbb{E}[(1+ch(X_2))^2] = 1 \times c^2$ . It appears that both inputs have the same contribution in the output variance when  $c = 1$ ,  $X_1$  is predominant when  $c < 1$  and  $X_2$  is predominant when  $c > 1$ . Moreover, when  $c$  tends to the infinity, the random function  $f(X)$  tends to  $ch(X_2)h(X_1)$  where the effects of  $X_1$  and  $X_2$  on the output are completely equal because of the symmetry of the model  $f$ .

### Conclusion about dependence measures vs. Sobol' indices

**These particular analyses reveal that the HSIC and the distance correlation nuance the conclusion obtained with the Sobol' indices. While the latter only focus on the input parameter contribution in the output variance, the dependence measures seem to be more sensitive to the global behavior of the output.** For the distance correlation, this can be explained by the fact that the distance covariance measures the distance between the product of the characteristic functions of a given input parameter and the model output and the characteristic function of the couple made by both variables. Consequently, it uses more information about the input-output relations because the characteristic function completely defines the probability distributions of these variables (separately and jointly). Furthermore, the HSIC maps the input and output values into the real line using some RKHS functions and measures the covariance between both functions; the associated estimator puts into relation the Gram matrices based on the associated reproducing kernels. In this way, the HSIC also uses more information about the output behavior than the Sobol' indices. **Finally, GSA conclusions can be strongly different between Sobol' and dependence measures in the presence of interaction effects.**

In the next section, the dependence measures are considered in a screening framework for high-dimensional problems, where the number of influential input parameters is (much) lower than the total ones. In these situations, industrial applications often require to eliminate the non-significant inputs before the computation of the Sobol' indices for the significant ones, which are sensitivity indices of great interest for engineers. Indeed, the Sobol' approach robustness requires a lot of model evaluations, especially for high-dimension problems, and are not tractable for the whole set of input parameters; a subset of relevant inputs must be selected. Moreover, we have noted in numerical experiments not mentioned in this paper that, in the presence of influential and non-influential input parameters,

the dependence measures take very different orders. Consequently, they arouse interest for screening purpose, in order to eliminate non-significant variables. It is in this sense we study the associated significance tests proposed and developed in Section 3.

## 5.2 Significance tests for screening

Now, we consider the function from Morris et al (2006) associating to the real input vector  $X = (X_1, \dots, X_d, X_{d+1}, \dots, X_{d+\check{d}})$  the scalar output

$$f(X) = a \sum_{i=1}^d \left( X_i + b \sum_{i < j=2}^d X_i X_j \right) \quad (6)$$

with  $a = \sqrt{12} - 6\sqrt{0.1(d-1)}$ ,  $b = 12\sqrt{0.1(d-1)}$  and  $X_i \stackrel{i.i.d.}{\sim} \mathcal{U}([0, 1])$ ,  $\forall i \in \{1, \dots, d + \check{d}\}$ . The  $d$  first input parameters are the influential inputs while the  $\check{d}$  are the non-influential ones. The ratio  $r = \frac{\check{d}}{d}$  is the quantity of non-significant variables brought back to the quantity of significant ones.

The objective of this section is to evaluate the potential of the different dependence measures in terms of screening using their associated significant tests presented and proposed in Section 3, for different sample sizes and different ratios  $r$ . A second objective is to study the screening performances of the Lasso regression and the bootstrap tests associated to the linear model (3) in Section 4.

### Comparison of different statistical tests based on dependence measures

For this model  $f$ , asymptotic and bootstrap tests based on the dependence measures mentioned in Section 2.1 are not at all satisfactory. Indeed, in the case of the Pearson's correlation coefficient with  $n = 500$  and  $\check{d} = d$  for example, the null hypothesis is kept for all input parameters, with a mean  $p$ -value equal to 1 for each influential factors, and to 0.5 for the others. In the same way, we obtain a mean  $p$ -value equal to 0.5 for the Spearman's correlation coefficient. This results are not surprising because the model  $f$  is not linear and so does not respect the hypotheses underlying to these coefficients. This analytical application illustrates the limitations of such correlation coefficients and justifies the use of the other dependence measures such as HSIC and distance correlation.

Table 3 compares the significance tests associated to the HSIC and distance correlation using 1000 Monte-Carlo runs for each pair  $(n, r)$  and computing the percentage of non-influential and influential input selection; the number of significant variables is equal to  $d = 5$ . Among these last quantities, the first one is the rate of the type I error and the second one is the power of the test<sup>2</sup>, usual notions in significance tests. Moreover, this table supplies the percentage of "perfect screening", which corresponds to the situation where all the non-significant variables are judged non-influential by the test while all the significant ones are judged influential. Lastly, the significance tests are presented in their asymptotical (Section 3.1), spectral (Section 3.2) and bootstrap (Section 3.3) versions for each dependence measure with a level  $c$  equal to 5%.

First of all, Table 3 shows that whatever the considered test, the rate of type I error and the power are independent of the non-significant input proportion  $r$ . Moreover, the rate of perfect screening increases with the number of observations  $n$  and decreases with  $r$ . It is also higher with the distance correlation tests. We also note that a powerful test does not imply an important perfect screening rate.

Then, considering the distance correlation, the results confirm the conservative property of the asymptotical test with a type I error around 1.5%; this implies a test power lower than using the asymptotical test based on HSIC. However, the power increases with  $n$  and this difference tends to disappear. Moreover, the bootstrap and the spectral tests give similar results, even if the spectral one is slightly conservative for a very small sample size while the bootstrap one is more powerful. Turning to HSIC, the asymptotical test has a type I error rate a little greater than the specified level (5%) when  $n$  is very small, while the spectral approach is slightly conservative.

Finally, we propose some advises to choose a statistical test according to the problem. A first point to note is the independence of the statistical tests from the proportion of non-significant variables. So, the remaining problem parameter lies in the number of observations:

- For a very small sample size (e.g.  $n < 50$  in this case), the bootstrap test for distance correlation is more powerful than the HSIC-based one.
- For a medium sample size (e.g.  $n \sim 100$ ), the distance correlation is also preferred to HSIC for the same reason, but this time, the spectral approach is advised for the significance test.

<sup>2</sup>The type I error occurs when the test concludes that a non-significant input is significant. The power of the test is the probability to conclude that a significant input is significant.

- When the number of observations is sufficiently important (e.g.  $n = 200$  for  $d + \check{d} = 55$  variables), all the tests agree on conclusions, except the asymptotical test based on distance covariance which has a better rate of perfect screening because of its conservative aspect. So, asymptotical tests are the better solutions when  $n$  is high, because of the previous reason for the distance covariance, and because of CPU time savings for both dependence measures.

Beyond the conclusions about the better approach (asymptotical, spectral or bootstrap), this comparison highlights that the tests based on distance covariance are often more powerful than those based on the HSIC. An explanation of this situation can be found in the definition of the distance covariance which measures the distance to the independence using characteristic functions, that is to say using the law definitions of the input parameters and of the output directly. However, the distance covariance is limited to vectorial inputs and outputs, contrarily to the HSIC which can deal with matricial inputs for example. Consequently, we advise to use the distance covariance for vectorial inputs and outputs and the HSIC for more complex data.

$n$	Significance test $\rightarrow$		Asymptotical			Spectral			Bootstrap		
		$r \rightarrow$	2	5	10	2	5	10	2	5	10
10	Non-influential	HSIC	7.6	7.3	7.5	3.7	3.8	3.8	5.0	5.0	5.0
		DCOR	<b>1.5</b>	<b>1.5</b>	<b>1.7</b>	4.4	4.2	4.3	4.6	4.9	4.8
	Influential	HSIC	<b>19.8</b>	<b>20.1</b>	<b>21.1</b>	13.1	14.1	13.6	15.3	16.3	16.2
		DCOR	10.7	11.2	11.4	<b>21.2</b>	<b>20.0</b>	<b>20.7</b>	<b>22.6</b>	<b>23.1</b>	<b>22.7</b>
	Perfect screening	HSIC	0.0	0.0	0.0	0.0	0.2	0.0	0.0	0.0	0.0
		DCOR	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
25	Non-influential	HSIC	6.8	3.1	5.8	5.7	4.9	4.6	5.9	5.4	5.0
		DCOR	<b>1.7</b>	<b>1.6</b>	<b>1.5</b>	5.0	4.5	4.9	5.4	5.1	4.7
	Influential	HSIC	<b>40.1</b>	<b>40.1</b>	<b>40.4</b>	37.7	37.0	37.5	38.5	38.6	38.7
		DCOR	37.9	37.4	37.7	<b>56.1</b>	<b>56.3</b>	<b>56.9</b>	<b>57.7</b>	<b>57.3</b>	<b>57.9</b>
	Perfect screening	HSIC	0.3	0.4	0.0	0.2	0.2	0.0	0.3	0.3	0.0
		DCOR	0.2	0.2	0.1	<b>3.0</b>	<b>2.0</b>	0.0	<b>2.6</b>	<b>2.6</b>	0.1
50	Non-influential	HSIC	5.3	5.2	5.3	4.6	4.6	4.7	5.0	4.8	4.9
		DCOR	<b>1.1</b>	<b>1.3</b>	<b>1.4</b>	4.7	4.6	5.0	4.6	4.5	4.7
	Influential	HSIC	70.8	70.0	71.5	69.2	68.4	70.1	69.9	68.8	70.4
		DCOR	<b>75.4</b>	<b>74.6</b>	<b>76.7</b>	<b>87.5</b>	<b>87.8</b>	<b>87.7</b>	<b>88.4</b>	<b>87.7</b>	<b>89.2</b>
	Perfect screening	HSIC	9.7	0.8	0.5	8.9	8.0	0.7	9.9	7.9	0.5
		DCOR	<b>19.3</b>	<b>17.2</b>	<b>9.4</b>	<b>35.6</b>	<b>29.6</b>	<b>4.0</b>	<b>41.4</b>	<b>30.1</b>	<b>4.7</b>
100	Non-influential	HSIC	5.1	5.3	5.4	4.9	4.9	5.1	4.9	4.8	5.0
		DCOR	<b>1.6</b>	<b>1.3</b>	<b>1.6</b>	4.4	4.7	5.1	4.9	4.7	5.0
	Influential	HSIC	95.7	95.9	95.9	95.5	95.8	95.7	95.6	95.7	95.8
		DCOR	<b>98.0</b>	<b>98.2</b>	<b>98.0</b>	<b>99.5</b>	<b>99.6</b>	<b>99.5</b>	<b>99.4</b>	<b>99.5</b>	<b>99.4</b>
	Perfect screening	HSIC	61.1	47.3	6.2	60.7	49.0	6.8	61.4	49.0	6.9
		DCOR	<b>83.9</b>	<b>80.8</b>	<b>40.7</b>	<b>77.9</b>	<b>60.2</b>	<b>9.1</b>	<b>75.5</b>	<b>60.1</b>	6.9
200	Non-influential	HSIC	4.7	5.3	5.2	4.5	5.0	4.9	4.5	4.9	5.0
		DCOR	<b>1.3</b>	<b>1.2</b>	<b>1.4</b>	4.2	5.2	5.1	4.7	5.0	5.0
	Influential	HSIC	99.9	99.9	100	99.9	99.9	100.0	99.9	99.9	100
		DCOR	100	100	100	100	100	100	100	100	100
	Perfect screening	HSIC	78.5	57.5	6.1	79.3	59.6	7.5	79.1	59.8	7.4
		DCOR	<b>93.4</b>	<b>88.5</b>	<b>47.5</b>	80.8	59.2	8.3	78.8	60.1	7.2

Table 3: Percentage of non-influential and influential input selection and perfect screening for different 5%-level significance tests, different sample sizes and different ratios of non-influential inputs, with HSIC and dCor. When a measure is significantly better than another, the percentage is in bold type. For each category (non-influential, influential and perfect screening), the best estimation approach is highlighted in gray.

From a computational point of view, we compare the different approaches for different sample size  $n$  and different ratio  $r$  of non-influential input parameters. We only consider the case of the HSIC because both kinds of sensitivity measures have similar formulations and, so, conclusions are similar. According to Figure 2, asymptotical approximation is clearly the cheapest approach while the bootstrap one (with  $B = 1000$  replications in this case, for a robustness purpose) is the most expensive. Moreover, contrarily to the spectral method, the resampling one has a computational

time very sensitive to  $n$  and  $r$ .

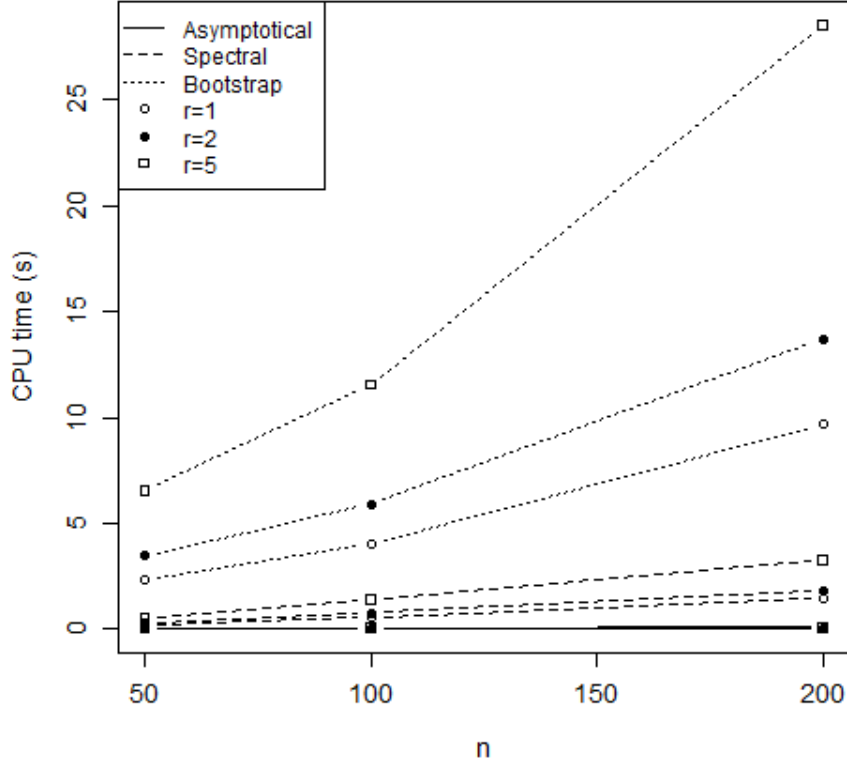


Figure 2: CPU time comparison for different sample sizes and different ratios of non-influential inputs with HSIC.

### Linear regression with HSIC

In this second part, we apply the linear model (3) based on HSIC to the numerical model (6) for  $d = \check{d} = 5$ . We consider the bootstrap significance tests proposed in Section 4.3 for the nested model selection and the HSIC Lasso with cross-validation proposed in Section 4.2. Table 4 gives the percentage of selected non-influential and influential inputs as well as the percentage of perfect screening, for different sample sizes and different methods: the bootstrap significance tests with 5%-level (bootstrap), the HSIC Lasso with cross-validation minimization (Lasso 1) and the HSIC Lasso with our improved cross-validation minimization (Lasso 2).  $N = 1000$  Monte-Carlo runs have been realized and, for the Lasso regression, we have adapted a Matlab implementation of the LARS algorithm<sup>3</sup> for the positive Lasso (Efron et al, 2004). Firstly, whatever the number of observations  $n$ , the bootstrap approach selects no more than 5% of non-influential inputs while the Lasso methods keep a lot of these variables: more than 60% with Lasso 1 and between 8 and 20% with Lasso 2, according to the sample size. On the contrary, the bootstrap approach is less powerful with small samples than the Lasso regression. Finally, the HSIC Lasso with our improved version of cross-validation minimization leads to a better perfect screening rate than using the classical one and when the ratio  $n/(d + \check{d})$  increases, the bootstrap approach is the method providing the more accurate screening.

To conclude:

- The bootstrap significance tests proposed in Section 4.3 for the nested model selection constitute the best approach for screening in a linear regression framework, except if  $n/(d + \check{d})$  is too small.
- Moreover, HSIC Lasso is an interesting tool but the choice of the penalization constant  $\lambda$  is an open-problem. Numerical tests reveal that its selection by cross-validation minimization leads to the selection of a too important number of non-influential variables. To cope with this issue, our improved cross-validation version seems to be a promising alternative.

<sup>3</sup>Matlab implementation of the LARS algorithm: <http://www.stat.berkeley.edu/~yugroup/downloads/>.

- Lastly, we remark that the bootstrap results are less good than those using the bootstrap distribution of the HSIC under the null hypothesis in Table 3.

$n$		Bootstrap	Lasso 1	Lasso 2
50	Non-influential	<b>4.0</b>	63.9	20.1
	Influential	68.2	<b>97.7</b>	85.7
	Perfect screening	7.4	1.4	<b>14.9</b>
100	Non-influential	<b>4.2</b>	66.3	14.1
	Influential	92.5	<b>99.9</b>	96.1
	Perfect screening	<b>48.6</b>	0.7	40.9
200	Non-influential	<b>4.5</b>	62.3	8.2
	Influential	99.8	100	99.6
	Perfect screening	<b>74.9</b>	2.6	66.6

Table 4: Percentage of non-influential and influential input selection and perfect screening for different feature selection approaches and different sample sizes, with HSIC in the linear regression model.

## 6 Nuclear reliability application

In the scope of nuclear reliability and nuclear power plant lifetime program, physical modeling tools have been developed to assess the component reliability of nuclear plants in numerous scenarios of use or accident. In the framework of nuclear plant risk assessment studies, the evaluation of component reliability during accidental conditions is a major issue required for safety studies. A workflow with chained thermal-hydraulic (TH) and thermal-mechanical (TM) computer codes models the behavior of the considered component subjected to highly hypothetical accidental conditions. It computes a safety criterion  $Y$  function of several uncertain parameters: 31 TH inputs and 10 TM inputs. Previous studies have fully characterized the probabilistic distributions of these 41 independent random variables and a sample of 400 independent realizations is available for a sensitivity analysis purpose.

The different significance tests based on the HSIC and presented in Section 3 are applied to this application, using the observations  $(X_1^{(i)}, \dots, X_{31}^{(i)}, \tilde{X}_1^{(i)}, \dots, \tilde{X}_{10}^{(i)}; Y^{(i)})_{1 \leq i \leq 1000}$ , where  $X_1, \dots, X_{31}$  are the TH inputs and  $\tilde{X}_1, \dots, \tilde{X}_{10}$  are the TM ones. The significance test level is fixed to  $\alpha = 5\%$ ; input parameters with a p-value greater than  $\alpha$  are ruled out and the remaining ones are considered as significant. Table 5 gives the p-values associated to the significant inputs for the different tests and supplies the HSIC values. Firstly, 50% of the output variability is due to the  $\tilde{X}_1$  uncertainty and 12% is due to the  $\tilde{X}_{10}$  uncertainty. The other influential variables have contributions to the output variability lower than 5%. Moreover, the input variables selected from a 5% significance test level have very small p-values and a level equal to 1% leads to the same selection. Consequently, the screening conclusions are highly probable. We have also increased the number of observations and have observed the robustness of these conclusions. The problem dimension, which is the number of model observations divided by the number of model inputs, is classical in presence of costly computer codes; nevertheless, we were able to add some tens of observations and saw the robustness of our results.

Finally, classical GSA tools such as Sobol' indices could be computed using thousands of simulations of a surrogate model built with few observations of the computer code and taking as input the inputs selected using the significance tests. In this way, these screening tools belong to a pretreatment phase conceived in order to increase the efficiency of the quantitative GSA.

	Inputs	$X_8$	$\tilde{X}_1$	$\tilde{X}_4$	$\tilde{X}_8$	$\tilde{X}_9$	$\tilde{X}_{10}$
1	100 $\times$ p-value	0.65	0.00	1.85	0.02	0.00	0.00
2	100 $\times$ p-value	1.40	0.00	2.80	0.00	0.00	0.00
3	100 $\times$ p-value	1.30	0.00	1.90	0.00	0.00	0.00
	100 $\times$ HSIC	0.15	2.71	0.13	0.19	0.25	0.65
	% w.r.t. the whole inputs	3%	50%	2%	3%	5%	12%
		3%			72%		
	% w.r.t. the influent inputs	4%	67%	3%	5%	6%	16%
		4%			96%		

Table 5: p-values (1 = estimator law replacement; 2 = spectral approximation; 3 = bootstrap) and HSIC for the nuclear application for the significant TH and TM inputs, with a significant test level equal to 5% and  $n = 400$ .



## 7 Conclusion

In this paper, we introduce recent developments around the use of dependence measures for sensitivity analysis (SA) and screening purposes. This situation occurs notably during the first steps of the establishment of a model, when the influential inputs are not exactly known and the precaution requires to consider all the potentially significant variables. Because of the costly nature of the numerical simulator, only some observations can be obtained, which prevents the use of classical SA quantitative methods, such as Sobol' indices, for these high-dimensional problems. Furthermore, classical Sobol' indices only focus on the decomposition of the output variance and not on its entire probabilistic distribution. For all these reasons, we turn to dependence measures recently proposed for global sensitivity analysis: the distance correlation and the Hilbert-Schmidt independence criterion (HSIC). The HSIC considers the covariance between two RKHS functions applied to these variables, and the distance covariance leading to the distance correlation corresponds to the mean norm between the characteristic function of both variables and the product of the characteristic functions of these variables.

At first, considering a sparse problem where the number of non-significant input parameters can be very important, independence hypothesis tests are required to use these new measures directly for a screening purpose. For this, asymptotic versions of such tests exist. Spectral approximations for the probabilistic laws involved in the asymptotic tests could improve some intrinsic approximations, especially in the presence of a medium size sample. From this, we propose non-asymptotic versions for these independence tests, in the case where the number of observations is low compared to the number of uncertain inputs. These non-asymptotic tests are based on a bootstrap sampling method. Always for a screening purpose, we propose a second approach based on the decomposition of any local measure of difference between two observed outputs as a linear regression on the same measures between the corresponding inputs. The regression coefficients are estimated using a linear least-squares minimization under positivity constraints. A coefficient equal to zero means that the corresponding input has no significant influence on the output. Thus, testing the nullity of each coefficient provides a screening method. In the case of the HSIC, we show that this model with  $\ell^1$ -penalization corresponds to the HSIC Lasso approach for feature selection and we propose to solve this problem using the LARS algorithm with positive coefficients. We also introduce a method for the selection of the penalty constant, based on the minimization of the cross-validation error reduced by a weighting of the associated standard-deviation. Likewise, we propose to apply the classic tools of model selection and, in particular, a bootstrap method testing the nullity of the model coefficients. To compare the different proposed approaches for screening based on dependence measures, we performed several numerical tests on classical analytical functions. Concerning the first approach, these experiments show that the different proposed significant tests based on dependence measures are very efficient. The ones based on distance correlation are sometimes more powerful while the ones based on HSIC have the advantage to be well-adapted to the case of high dimensional inputs. Concerning the kind of significance test (asymptotic, spectral and non-asymptotic), the compromise "CPU time - accuracy" gives the advantage to the bootstrap tests in the presence of small sample sizes and to the asymptotical approaches when the number of observations is higher. The spectral approximation of the asymptotical law can be viewed as an intermediary solution between these two extreme configurations. In addition, the first approach using directly the dependence measures seems to be slightly better than the second one, based on the linear decomposition of these sensitivity measures. Some of these methods are successfully applied to a nuclear reliability application.

In this paper, we also try to provide some preliminary answers to the question "What sensitivity indices to what situation?", without pretension to build a theory. For this, we performed many tests on toy functions to compare the results given by the HSIC, the distance correlation and the Sobol' index. Firstly, the new dependence measures lead to conclusions of the same order than the Sobol' indices ones. Then, they seem to be higher for the linear effects than for the non-linear ones, these effects being additive, centered and with variance equal to one, which leads to uniform Sobol' indices. Moreover, the HSIC further detects the monotonic effects v.s. the non-monotonic ones while the opposite occurs with the distance correlation. These tests also highlights that the dependence measures are more sensitive to the presence of an interaction term than the Sobol' indices and yield some different sensibility analysis conclusions. Their interpretation seems closer to that of the density-based sensitivity indices such as Borgonovo's. Beyond this complementary aspect, the HSIC and the distance correlation need only a few number of model evaluations, which is a great advantage over the classical variance-based or density-based indices. Finally, various numerical tests illustrate that dependence measures provide a relevant information which is coherent and sometimes complementary to the one obtained with classical indices.

Given the above, we advise the use of dependence measures associated to independence tests in global sensitivity analysis when the number of simulations is weak, when the problem takes place in an high-dimensional context or when we want to reinforce or qualify the conclusions obtained with the Sobol' indices. Moreover for industrial problems, the aim of the GSA is often to reduce the output variance of the simulator. In this case, the use of dependence measures can be viewed as a selection step and then, a quantitative phase consists to compute the Sobol' indices of the retained model inputs. Da Veiga (2014) also shows the interest of such sensitivity measures

for high-dimensional output. However, in the presence of many thousands of observations, distance correlation and HSIC estimations are CPU time-expensive and other kernel methods should be investigated. More recently, a new dependence measure called “randomized dependence coefficient” has been proposed (López-Paz et al, 2013) with a computational cost of  $\mathcal{O}(n \log(n))$  while the distance correlation and HSIC ones are of  $\mathcal{O}(n^2)$ ,  $n$  being the number of required simulations. Considering this coefficient for GSA problems could be an interesting extension to this paper. In addition, applying the present screening methods to industrial applications, with functional inputs and outputs, could be a follow-up to this work. Finally, it should be interesting to study the dependence of the significance test results to the kernel functions chosen to compute the HSIC.

## Acknowledgments

We are grateful to Béatrice Laurent and Sébastien Da Veiga for helpful discussions. We thank Simon Nanty for his technical help in the nuclear application.

## References

- Santner TJ, Williams BJ, Notz W (2003) The design and analysis of computer experiments. Springer series in statistics, Springer, New York
- de Rocquigny E, Devictor N, Tarantola S (2008) Uncertainty in industrial practice. Wiley
- Saltelli A, Ratto M, Andres T, Campolongo F, Cariboni J, Gatelli D, Saisana M, Tarantola S (2008) Global sensitivity analysis: the primer. John Wiley
- Kleijnen JP (1997) Sensitivity analysis and related analyses: A review of some statistical techniques. *Journal of Statistical Computation and Simulation* 57(1-4):111–142
- Frey HC, Patil SR (2002) Identification and Review of Sensitivity Analysis Methods. *Risk Analysis* 22(3):553–578
- Helton J, Johnson J, Sallaberry C, Storlie C (2006) Survey of sampling-based methods for uncertainty and sensitivity analysis. *Reliability Engineering & System Safety* 91(10 – 11):1175 – 1209
- Cacuci DG (1981) Sensitivity theory for nonlinear systems. I. Nonlinear functional analysis approach. *Journal of Mathematical Physics* 22(12):2794–2802
- Lamboni M, Iooss B, Popelin AL, Gamboa F (2013) Derivative-based global sensitivity measures: General links with sobol’ indices and numerical tests. *Mathematics and Computers in Simulation* 87:45–54
- Sobol I (1993) Sensitivity estimates for nonlinear mathematical models. *MMCE* 1:407–414
- Borgonovo E (2007) A new uncertainty importance measure. *Reliability Engineering & System Safety* 92(6):771 – 784
- Saltelli A, Annoni P (2010) How to avoid a perfunctory sensitivity analysis. *Environmental Modelling & Software* 25(12):1508 – 1517
- Daniel C (1958) On varying one factor at a time. *Biometrics* 14:430–431
- Morris MD (1991) Factorial sampling plans for preliminary computational experiments. *Technometrics* 33(2):161–174
- Bettonvil B, Kleijnen JP (1997) Searching for important factors in simulation models with many factors: Sequential bifurcation. *European Journal of Operational Research* 96(1):180–194
- Montgomery DC (2006) Design and Analysis of Experiments. John Wiley & Sons
- De Castro Y, Janon A (2014) Randomized pick-freeze for sparse Sobol indices estimation in high dimension. Research report, <http://hal.inria.fr/hal-00962473>, Laboratoire de Mathématiques d’Orsay
- Caflich R, Morokoff W, Owen A (1997) Valuation of mortgage backed securities using brownian bridges to reduce effective dimension. *Journal of Computational Finance* 1 pp 27–46
- Marrel A, Iooss B, Laurent B, Roustant O (2009) Calculations of sobol indices for the gaussian process metamodel. *Rel Eng & Sys Safety* 94(3):742–751
- Sudret B (2008) Global sensitivity analysis using polynomial chaos expansions. *Reliability Engineering & System Safety* 93(7):964 – 979

- Gretton A, Bousquet O, Smola A, Schölkopf B (2005) Measuring statistical dependence with hilbert-schmidt norms. In: *Proceedings Algorithmic Learning Theory*, Springer-Verlag, pp 63–77
- Székely GJ, Rizzo ML, Bakirov NK (2007) Measuring and testing dependence by correlation of distances. *The Annals of Statistics* 35(6):2769–2794
- Blaschko MB, Zaremba W, Gretton A (2013) Taxonomic prediction with tree-structured covariances. In: *Machine Learning and Knowledge Discovery in Databases, Lecture Notes in Computer Science*, vol 8189, Springer Berlin Heidelberg, pp 304–319
- Da Veiga S (2014) Global sensitivity analysis with dependence measures. *Journal of Statistical Computation and Simulation*
- Kendall M, Stuart A (1977) *The advanced theory of statistics*, vol 2: Inference and Relationship, 4th edn. Charles Griffin, London
- Aronszajn N (1950) Theory of Reproducing Kernels. *Transactions of the American Mathematical Society* 68(3)
- Deza M, Deza E (2009) Encyclopedia of distances. In: *Encyclopedia of Distances*, Springer Berlin Heidelberg, pp 1–583
- Fukumizu K, Gretton A, Lanckriet GR, Schölkopf B, Sriperumbudur BK (2009) Kernel choice and classifiability for rkhs embeddings of probability distributions. In: Bengio Y, Schuurmans D, Lafferty J, Williams C, Culotta A (eds) *Advances in Neural Information Processing Systems 22*, Curran Associates, Inc., pp 1750–1758
- Balasubramanian K, Sriperumbudur BK, Lebanon G (2013) Ultrahigh dimensional feature screening via rkhs embeddings. In: *AISTATS, JMLR Proceedings*, vol 31, pp 126–134
- Yamada M, Jitkrittum W, Sigal L, Xing EP, Sugiyama M (2014) High-dimensional feature selection by feature-wise kernelized lasso. *Neural Computation* 26(1):185–207
- Gretton A, Fukumizu K, Teo CH, Song L, Schölkopf B, Smola AJ (2007) A kernel statistical test of independence. In: Platt JC, Koller D, Singer Y, Roweis ST (eds) *NIPS*, Curran Associates, Inc.
- Shen H, Jegelka S, Gretton A (2009) Fast kernel-based independent component analysis. *IEEE Transactions on Signal Processing* 57(9):3498–3511
- Sejdinovic D, Sriperumbudur B, Gretton A, Fukumizu K (2013) Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics* 41(5):2263–2291
- Zhang K, Peters J, Janzing D, Schölkopf B (2011) Kernel-based conditional independence test and application in causal discovery. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)* 804–813. AUAI Press, Corvallis, Oregon.
- Peng H, Long F, Ding C (2005) Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27:1226–1238
- Hoerl AE, Kennard RW (1970) Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* 12:55–67
- Tibshirani R (1996) Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society (Series B)* 58:267–288
- Zou H, Hastie T (2005) Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B* 67:301–320
- Efron B, Hastie T, Johnstone I, Tibshirani R (2004) Least angle regression. *Annals of Statistics* 32:407–499
- Kleijnen J (2007) *Design and Analysis of Simulation Experiments*. International Series in Operations Research & Management Science, Springer
- Plischke E, Borgonovo E, Smith CL (2013) Global sensitivity measures from given data. *European Journal of Operational Research* 226(3):536 – 550
- López-Paz D, Hennig P, Schölkopf B (2013) The randomized dependence coefficient. In: *Advances in Neural Information Processing Systems 26*, pp 1–9
- Morris MD, Moore LM, McKay MD (2006) Sampling plans based on balanced incomplete block designs for evaluating the importance of computer model inputs. *Journal of Statistical Planning and Inference* 136(9):3203 – 3220